

Estatística Inferencial

Prof. Wagner Hugo Bonat

Departamento de Estatística
Universidade Federal do Paraná





Verossimilhança e suas derivadas

Notação e definições (relembrando)

- ▶ $Y = (Y_1, \dots, Y_n)^T$: v.a.'s independentes e idênticamente distribuídas.
- ▶ $Y_i \sim f(\theta)$ onde f denota a função densidade de probabilidade ou função de probabilidade e $\theta = (\theta_1, \dots, \theta_p)^T$ é um vetor de p parâmetros populacionais.
- ▶ $y = (y_1, \dots, y_n)^T$ denota o vetor de valores observados da v.a. Y .
- ▶ **Estatística** - Uma estatística é uma variável aleatória $T = t(Y)$, onde a função $t(\cdot)$ não depende de θ .
- ▶ **Estimador** - Uma estatística T é um estimador para θ se o valor realizado $t = t(y)$ é usado como uma estimativa para o valor de θ .
- ▶ **Distribuição amostral** - A distribuição de probabilidade de T é chamada de distribuição amostral do estimador $t(Y)$.

Função de verossimilhança (caso uniparamétrico)

- ▶ **Função de verossimilhança** - Seja \mathbf{y} um vetor $n \times 1$ representando uma realização de um vetor aleatório \mathbf{Y} com função de probabilidade ou densidade probabilidade $f(\mathbf{Y}, \theta)$, onde θ denota um parâmetro, com $\theta \in \Theta$, sendo Θ o respectivo espaço paramétrico. A função de verossimilhança ou simplesmente verossimilhança para θ dado os valores observados \mathbf{y} é a função aleatória $L(\theta|\mathbf{y}) \equiv f(\mathbf{Y}, \theta)$.

- ▶ Caso discreto:

$$L(\theta|\mathbf{y}) \equiv P_{\theta}(\mathbf{Y} = \mathbf{y}).$$

- ▶ Caso contínuo (simplificado)

$$L(\theta|\mathbf{y}) \approx \prod_{i=1}^n f(y_i, \theta).$$

Verossimilhança - Condições de regularidade

- ▶ O parâmetro θ é **identificável**. Isso significa que se $f(\theta_1|\mathbf{y}) = f(\theta_2|\mathbf{y})$ para quase todos $\mathbf{y} \in \mathbb{R}$, então $\theta_1 = \theta_2$.
- ▶ Para quase todo significa que a condição não é verdadeira para um conjunto de \mathbf{y} com probabilidade zero de ocorrência.
- ▶ O suporte de $f(\theta|\mathbf{y})$ é o mesmo para todo $\theta \in \mathbb{R}$.
- ▶ O verdadeiro valor do parâmetro θ_0 pertence ao interior de Θ .
- ▶ $f(\theta|\mathbf{y})$ é duas vezes continuamente diferenciável com relação θ para quase todo $\mathbf{y} \in \mathbb{R}$.
- ▶ $\frac{\partial}{\partial \theta}$ e \int (caso contínuo) ou $\frac{\partial}{\partial \theta}$ e \sum (caso discreto) podem ser intercambiada.

Log-Verossimilhança (caso uniparamétrico)

- ▶ A função de log-verossimilhança é a função estocástica $l(\theta) : \Theta \rightarrow \mathbb{R}$ definida por

$$l(\theta|\mathbf{y}) = \log (L(\theta|\mathbf{y})) .$$

- ▶ No caso iid, tem-se

$$l(\theta|\mathbf{y}) = \sum_{i=1}^n \log (L(\theta|y_i)) .$$

- ▶ $l(\theta|\mathbf{y}) = -\infty$ quando $L(\theta) = 0$, mas isso ocorre quando $f(y_1, \dots, y_n|\theta) = 0$ que tem probabilidade de ocorrência igual a zero.

Desigualdade de Jensen

- ▶ **Desigualdade de Jensen:** Seja g uma função estritamente convexa e Y uma v.a com $E(|Y|) < \infty$ tal que a distribuição de Y é não degenerada. Então,
 - ▶ $g(E(Y)) < E(g(Y))$. Por outro lado, se g é estritamente côncava, então
 - ▶ $g(E(Y)) > E(g(Y))$.
- ▶ Lembrete:
 - ▶ Função convexa $f(tx + (1 - t)y) \leq tf(x) + (1 - t)f(y)$.
 - ▶ Função côncava $f(tx + (1 - t)y) \geq tf(x) + (1 - t)f(y)$.

Desigualdade de Jensen e Máxima Verossimilhança

- ▶ Teorema: Seja θ_0 o verdadeiro valor do parâmetro. Então,

$$P_{\theta_0}(L(\theta_0|\mathbf{y}) > L(\theta|\mathbf{y})) \rightarrow 1, \text{ quando } n \rightarrow \infty.$$

- ▶ Interpretação: $L(\theta_0) > L(\theta)$ com alta probabilidade para n grande.
- ▶ Assim, $L(\theta)$ vai tender a ter o seu máximo próximo a θ_0 , o verdadeiro valor de θ .
- ▶ Motiva a ideia de estimação por máxima verossimilhança.
- ▶ Demonstração (vídeo separado opcional).

Função escore e Informação de Fisher (caso uniparamétrico)

- Função escore para θ (efficient score)

$$\begin{aligned}U(\theta|Y) &= U(\theta|Y_1, \dots, Y_n) \\&= \sum_{i=1}^n \frac{\partial}{\partial \theta} l(\theta, Y_i). \\&= \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(\theta, Y_i).\end{aligned}$$

- Informação de Fisher ou Informação esperada

$$I_E(\theta) = \text{Var}(U(\theta|Y)).$$

- Informação de Fisher também é chamada de *intrinsic accuracy*.

Igualdades de Bartlett (caso uniparamétrico)

- ▶ Sob condições de regularidade, tem-se
 - ▶ Primeira igualdade: $E(U(\theta|Y)) = 0$.
 - ▶ Segunda igualdade: $I_E(\theta) = -E(l''(\theta|Y)) = -E(U'(\theta|Y))$.
- ▶ Implicação: $\text{Var}(U(\theta|Y)) = E(U(\theta|Y)^2)$.
- ▶ Demonstração.
- ▶ Exercício: Sejam $Y_i \sim P(\lambda)$ iid, para $i = 1, \dots, n$. Verifique as igualdades de Bartlett.

Informação observada

- Informação observada para θ

$$I_O(\theta) = -l''(\theta|Y).$$

- Note que $I_E(\theta) = E(I_O(\theta))$.
- Além disso, pela lei dos grandes números

$$I_O(\theta) \xrightarrow{P} I_E(\theta) \quad \text{quando } n \rightarrow \infty.$$

- Exercício: Sejam $Y_i \sim B(p)$ iid, para $i = 1, \dots, n$. Encontre a informação observada e esperada e mostre que a informação esperada coincide com a variância da função escore.

Desigualdade de Cramér-Rao (caso uniparamétrico)

- ▶ Teorema: Se $T(Y_1, \dots, Y_n)$ é um estimador não viciado para θ , então

$$\text{Var}(T) \geq I_E(\theta)^{-1}.$$

- ▶ A quantidade $I_E(\theta)^{-1}$ é chamado de limite inferior de Cramér-Rao.
- ▶ Um estimador não viciado é chamado eficiente se $\text{Var}(T) = I_E(\theta)^{-1}$.
- ▶ Demonstração.



Vetor de parâmetros

Função de verossimilhança (caso multiparamétrico)

- ▶ Função de verossimilhança $L : \Theta \rightarrow [0, \infty]$ é uma função aleatória de um **vetor** de argumentos definida por

$$L(\theta|\mathbf{y}) = \prod_{i=1}^n f(y_i, \theta), \quad \text{para } \theta \in \Theta.$$

- ▶ Função de log-verossimilhança $l : \Theta \rightarrow \mathbb{R}$ é uma função aleatória de um **vetor** de argumentos definida por

$$l(\theta|\mathbf{y}) = \sum_{i=1}^n \log f(y_i, \theta).$$

- ▶ Equivalentemente,

$$l(\theta|\mathbf{y}) = \log L(\theta|\mathbf{y}).$$

Vetor escore

- ▶ O vetor escore $U(\theta|\mathbf{y}) : \Theta \rightarrow \mathbb{R}^p$ é um vetor aleatório $p \times 1$ definido por

$$U(\theta|\mathbf{y}) = \frac{\partial l(\theta|\mathbf{y})}{\partial \theta} = \begin{pmatrix} \frac{\partial l(\theta|\mathbf{y})}{\partial \theta_1} \\ \vdots \\ \frac{\partial l(\theta|\mathbf{y})}{\partial \theta_p} \end{pmatrix}.$$

- ▶ Notação popular em termos de gradiente

$$U(\theta|\mathbf{y}) = \nabla_{\theta} l(\theta|\mathbf{y}).$$

- ▶ Notação j -ésimo componente de $U(\theta|\mathbf{y})$ por $U_j(\theta|\mathbf{y})$.

Esperança do vetor escore

- O vetor escore satisfaz as igualdades de Bartlett, ou seja,

$$E(U(\theta|Y)) = 0,$$

isso significa que

$$E(U_j(\theta|Y)) = 0, \quad \text{para } j = 1, \dots, p.$$

Matriz de informação esperada

- ▶ A matriz $p \times p$ definida por

$$\begin{aligned} I(\boldsymbol{\theta}) &= \text{Var}(U(\boldsymbol{\theta}|\mathbf{Y})) \\ &= E(U(\boldsymbol{\theta}|\mathbf{Y})U^T(\boldsymbol{\theta}|\mathbf{Y})). \end{aligned}$$

é chamada de matriz de informação esperada.

- ▶ As entradas j e k são expressadas por

$$I_{jk}(\boldsymbol{\theta}) = \text{Cov}(U_j(\boldsymbol{\theta}|\mathbf{Y}), U_k(\boldsymbol{\theta}|\mathbf{Y})) \quad (1)$$

$$= E(U_j(\boldsymbol{\theta}|\mathbf{Y}), U_k(\boldsymbol{\theta}|\mathbf{Y})). \quad (2)$$

Matriz de informação observada

- ▶ A matriz $p \times p$ definida por

$$J(\theta) = -\frac{\partial^2 l(\theta|\mathbf{y})}{\partial \theta \partial \theta^\top}.$$

- ▶ As entradas j e k da matriz de informação observada é dada por

$$J_{jk}(\theta) = -\frac{\partial^2 l(\theta|\mathbf{y})}{\partial \theta_j \partial \theta_k}.$$

- ▶ Segunda igualdade de Bartlett

$$I(\theta) = E(J(\theta)).$$

Desigualdade de Cramér-Rao generalizada

- ▶ Defina $I^{jk}(\boldsymbol{\theta}) = \{I^{-1}(\boldsymbol{\theta})\}_{jk}$. Se $T = T(Y_1, \dots, Y_n)$ é um estimador não-viciado para θ_1 , ou seja,

$$E(T) = \theta_1,$$

então

$$\text{Var}(T) \geq I^{11}(\boldsymbol{\theta}).$$

- ▶ Demonstração análoga ao caso univariado usando a desigualdade de Cauchy-Schwarz generalizada.

Parâmetros ortogonais

- Considere um modelo estatístico parametrizado por $\theta = (\theta_1, \theta_2)^\top$. No caso da matriz de informação de Fisher ser diagonal

$$I(\theta) = \begin{pmatrix} I_{11}(\theta) & 0 \\ 0 & I_{22}(\theta) \end{pmatrix},$$

os parâmetros θ_1 e θ_2 são ditos **ortogonais**.

- O inverso da informação de Fisher é também diagonal

$$I^{-1}(\theta) = \begin{pmatrix} 1/I_{11}(\theta) & 0 \\ 0 & 1/I_{22}(\theta) \end{pmatrix}.$$

- Assim, $\hat{\theta}_1$ e $\hat{\theta}_2$ são assintoticamente independentes. com distribuição

$$\hat{\theta}_j \stackrel{a}{\sim} N(\theta_j, 1/I_{jj}(\theta)).$$

Parâmetros ortogonais: Generalização

- ▶ Considere um modelo estatístico parametrizado por $\theta = (\theta_1, \theta_2)^\top$. No caso da matriz de informação de Fisher ser bloco diagonal

$$I(\theta) = \begin{pmatrix} I_1(\theta) & 0 \\ 0 & I_2(\theta) \end{pmatrix},$$

os vetores de parâmetros θ_1 e θ_2 são ditos **ortogonais**.

- ▶ A distribuição assintótica de θ_1 é a mesma se θ_2 é considerado conhecido ou desconhecido.
- ▶ Definição similar pode ser feita usando a matriz de informação observada.

Exercício

- ▶ Sejam $Y_i \sim N(\mu, \sigma^2)$ iid para $i = 1, \dots, n$.
 - ▶ Escreva a função de verossimilhança e log-verossimilhança.
 - ▶ Obtenha a função escore.
 - ▶ Obtenha a matriz de informação observada e esperada.
 - ▶ Obtenha a variância assintótica de $\hat{\mu}$ e $\hat{\sigma}^2$.
 - ▶ Os parâmetros μ e σ^2 são ortogonais? Justifique.